

Faculty of Engineering and Information Technology  
University of Technology Sydney

# **Data Mining In Epigenetic Modification and Gene Expression**

A thesis submitted in partial fulfillment of  
the requirements for the degree of  
**Doctor of Philosophy**

by

Zhixun Zhao

October 2020



## **CERTIFICATE OF AUTHORSHIP/ORIGINALITY**

I, Zhixun Zhao declare that this thesis, is submitted in fulfilment of the requirements for the award of Doctor of Philosophy, in the Faculty of Engineering and Information Technology at the University of Technology Sydney.

This thesis is wholly my own work unless otherwise referenced or acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

This document has not been submitted for qualifications at any other academic institution.

This research is supported by the Australian Government Research Training Program.

Production Note:

SIGNATURE: Signature removed prior to publication.  
[Zhixun Zhao]

DATE: 15<sup>th</sup> October, 2020

PLACE: Sydney, Australia



# Acknowledgments

First and foremost, I want to thank my supervisor, Prof. Jinyan Li, for his extensive instruction and patient guidance during the last three and a half years, regardless of my research or daily life. His guidance helped me improve my research skills, such as scientific writing, academic communication, and presentation. I appreciate all his contributions of time, ideas, and funding to make all of this thesis possible.

I want to thank my advisor, Prof. Liang Fang, who advised me at my home university in China. His strong support and help gave me the chance to win the CSC (China scholarship council) scholarship and study overseas. I thank Prof. Shaoqing Li and Prof. Jihua Chen, who supervised me when I was pursuing my master's degree, for their initial guidance on my research. I thank my co-supervisor, Prof. Fang Chen, for her help in my research and providing me the opportunity to participate in industry projects.

My sincere thanks also go to my research team members: Dr. Hui peng, Dr. Yi Zheng, Dr. Yuansheng Liu, Dr. Chaowang Lan, Xiaocai Zhang, Xuan Zhang, Tao Tang, and Tian Lan, for their help in both my research and life. It's my honor to be one member of my research team, which is more like a family. Thank you all for bringing me the feeling of home and the unforgettable memories. I have special thanks to Dr. Huijun Wu, who is a truly loyal friend to me. We live together for two years and many thanks for your company.

I am also grateful to acknowledge the funding sources, including the tuition fee and living expenses provided by the China Scholarship Council

## *Acknowledgments*

---

and Graduate Research School, travel funds provided by the Faculty of Engineering and Information Technologies, and vice-chancellor funding. Thanks to all staff of the Advanced Analytics Institute and School of Computer Science who provide services and conveniences to my study and research in UTS.

Lastly, I owe special thanks to my wife Zhujun Xue, my parents, and my parents-in-law. It's your love and encouragement that support me in completing my Ph.D. study. I dedicate this work to you all. Especially, I want to thank my wife for her understanding and sacrifice all these years. Love you forever.

Zhixun Zhao

October 2020 @ UTS

# Contents

|  |                       |
|--|-----------------------|
| Certificate . . . . .  | <b>i</b>              |
| Acknowledgment . . . . .   | <b>iii</b>            |
| List of Figures . . . . .  | <b>ix</b>             |
| List of Tables . . . . .   | <b>xi</b>             |
| List of Publications . . . . .   | <b>xiii</b>           |
| Abstract . . . . .   | <b>xv</b>             |
| <br><b>Chapter 1 Introduction . . . . .</b>  | <br><b>1</b>          |
| 1.1 Background . . . . .   | <b>1</b>              |
| 1.1.1 Epigenetic modification . . . . .  | <b>1</b>              |
| 1.1.2 Gene expression biomarker . . . . .  | <b>4</b>              |
| 1.1.3 Data mining in bioinformatics . . . . .  | <b>6</b>              |
| 1.2 Research questions . . . . .   | <b>8</b>              |
| 1.3 Research contributions . . . . .   | <b>12</b>             |
| 1.4 Thesis structure . . . . .   | <b>14</b>             |
| <br><b>Chapter 2 Related work and literature review . . . . .</b>  | <br><b>16</b>         |
| 2.1 DNA N <sup>4</sup> -methylcytosine prediction . . . . .  | <b>16</b>             |
| 2.2 mRNA N <sup>6</sup> -methyladenosine prediction . . . . .  | <b>18</b>             |
| 2.3 Lung cancer gene markers identification . . . . .  | <b>21</b>             |
| 2.4 Summary . . . . .  | <b>23</b>             |
| <br><b>Chapter 3 Accurate prediction of DNA N<sup>4</sup>-methylcytosine<br/>          sites via boost-learning various types of sequence<br/>          features . . . . .</b> | <br><br><br><b>25</b> |

|  |   |    |
|--|---|----|
| 3.1  | Background . . . . .  | 25 |
| 3.2  | Materials and methods . . . . .   | 27 |
| 3.2.1  | Benchmark datasets . . . . .  | 27 |
| 3.2.2  | Feature space construction . . . . .  | 28 |
| 3.2.3  | Feature selection scheme . . . . .  | 32 |
| 3.2.4  | Support vector machine . . . . .  | 33 |
| 3.2.5  | Performance evaluation metrics . . . . .  | 33 |
| 3.3  | Results . . . . .   | 34 |
| 3.3.1  | Feature importance analysis . . . . .   | 34 |
| 3.3.2  | Impact of feature selection on classification . . . . .                                   | 36 |
| 3.3.3  | Comparison with state-of-art predictors . . . . .   | 37 |
| 3.3.4  | Case study . . . . .  | 40 |
| 3.4  | Discussion and summary . . . . .  | 42 |
| <br><b>Chapter 4 Imbalance learning for the prediction of N<sup>6</sup>-methylation sites in mRNAs . . . . .</b> |   |    |
| 4.1  | Background . . . . .  | 44 |
| 4.2  | Materials and methods . . . . .   | 46 |
| 4.2.1  | Feature space construction . . . . .  | 47 |
| 4.2.2  | Imbalance learning . . . . .  | 51 |
| 4.2.3  | Performance evaluation metrics . . . . .  | 52 |
| 4.3  | Results . . . . .   | 53 |
| 4.3.1  | Specific SNP status as new features . . . . .   | 53 |
| 4.3.2  | Performance on the independent dataset . . . . .  | 55 |
| 4.3.3  | Robust performance when tested on datasets with<br>different imbalance ratios . . . . .   | 55 |
| 4.3.4  | Performance on 1226 individual transcripts . . . . .                                      | 56 |
| 4.3.5  | Feature importance analysis . . . . .   | 57 |
| 4.4  | Case studies . . . . .  | 59 |
| 4.4.1  | m <sup>6</sup> A site prediction for c-Jun transcript . . . . .                           | 59 |
| 4.4.2  | m <sup>6</sup> A site prediction for a transcript related to HIV-1<br>infection . . . . . | 60 |



|  |  |    |
|--|--|----|
| 4.5  | Discussion . . . . .   | 62 |
| 4.6  | Summary . . . . .  | 62 |
| <br><b>Chapter 5 Identification of lung cancer gene markers through<br/>kernel maximum mean discrepancy and information<br/>entropy . . . . . 63</b> |  |    |
| 5.1  | Background . . . . .   | 63 |
| 5.2  | Materials and methods . . . . .  | 65 |
| 5.2.1  | Dataset . . . . .  | 65 |
| 5.2.2  | Gene marker identification framework . . . . .                               | 65 |
| 5.2.3  | Kernel maximum mean discrepancy . . . . .                                    | 66 |
| 5.2.4  | Boundary discovery method . . . . .  | 68 |
| 5.2.5  | GO and KEGG enrichment analysis . . . . .                                    | 69 |
| 5.2.6  | Conventional DEA method and machine learning evaluation<br>metrics . . . . . | 70 |
| 5.3  | Results . . . . .  | 70 |
| 5.3.1  | Gene differential expression between different tissue types                  | 71 |
| 5.3.2  | Identify marker genes in cancer development . . . . .                        | 72 |
| 5.3.3  | GO and KEGG pathway enrichment . . . . .                                     | 74 |
| 5.3.4  | Expression boundary identification . . . . .                                 | 76 |
| 5.4  | Discussion . . . . .   | 77 |
| 5.5  | Summary . . . . .  | 78 |
| <br><b>Chapter 6 Conclusions and future work . . . . . 79</b>  |  |    |
| 6.1  | Conclusions . . . . .  | 79 |
| 6.2  | Future work . . . . .  | 81 |
| <br><b>Chapter A Appendix: Methodology foundation . . . . . 84</b>   |  |    |
| A.1  | Applied statistical methods . . . . .  | 84 |
| A.1.1  | Information entropy . . . . .  | 84 |
| A.1.2  | Fisher's exact test . . . . .  | 85 |
| A.2  | Adopted machine learning algorithms . . . . .                                | 85 |

|                     |   |           |
|---------------------|---|-----------|
| A.2.1               | Support vector machine . . . . .                  | 85        |
| A.2.2               | XGBoost . . . . .                                 | 86        |
| A.3                 | Cross validation and evaluation metrics . . . . . | 87        |
| A.3.1               | Cross validation . . . . .                        | 87        |
| A.3.2               | Performance evaluation metrics . . . . .          | 87        |
| <b>Chapter B</b>    | <b>Additional files . . . . .</b>                 | <b>89</b> |
| <b>Chapter C</b>    | <b>Appendix: List of Symbols . . . . .</b>        | <b>90</b> |
| <b>Bibliography</b> | <b>. . . . .</b>                                  | <b>92</b> |

# List of Figures

|     |   |    |
|-----|---|----|
| 1.1 | <b>Thesis structure.</b> It consists of the following four parts: introduction, related work, my work, and conclusions and future work. Short introduction of each part is shown in the right side. . . . .   | 15 |
| 3.1 | <b>Framework of proposed model construction</b> . . . . .   | 26 |
| 3.2 | <b>Sequence logos for DNA samples in the benchmark datasets</b> . . . . .   | 29 |
| 3.3 | <b>Sequence feature importance distribution</b> . . . . .   | 35 |
| 3.4 | <b>The independent test performanc before and after feature selection</b> . . . . .   | 37 |
| 3.5 | <b>The confidence of predicted label in case studies</b> . . . .  | 41 |
| 4.1 | <b>Feature space construction</b> . . . . .   | 47 |
| 4.2 | <b>SNP specificity ranking</b> The black blocks stand for the Fisher’s exact test rankings and the green blocks stand for the MRMR rankings. X-axis is the window sequence sites from -25 to 25. Y-axis is the total ranking of each position. A low ranking means a high SNP specificity at this position. . . . . | 54 |
| 4.3 | <b>Performance on datasets of different imbalance levels</b> The F1 and MCC values of four predictors are represented. X-axis k is the ratio of the negative samples to positive samples (imbalance level) in a test dataset; Y-axis is metric value. . . .   | 56 |
| 4.4 | <b>Boxplot of feature importance scores</b> . . . . .   | 58 |

|     |  |   |    |
|-----|--|---|----|
| 4.5 | <b>Predicted m<sup>6</sup>A sites in the case studies</b>        | The x-axis stands for the potential m <sup>6</sup> A sites confirming to the sequence motif DRACH and the y-axis indicates the four predictors. All colored blocks are the predicted m <sup>6</sup> A sites. Red blocks represent true positive sites, and yellow blocks are false positive ones. (a) the prediction results for the c-Jun case and (b) the predictions for the HIV-1 case. . . . . | 60 |
| 5.1 | <b>Gene marker identification framework</b>                      | . . . . .   | 66 |
| 5.2 | <b>Box-plot of gene expression levels in three tissue types.</b> | The X-axis is the FPKM expression level; the Y-axis is the tissue type. . . . .   | 74 |
| 5.3 | <b>KEGG pathway enrichment analysis for top ranking genes.</b>   | . . . . .   | 76 |

# List of Tables

|     |   |           |
|-----|---|-----------|
| 3.1 | <b>Summary of six benchmark datasets</b> . . . . .  | <b>28</b> |
| 3.2 | <b>The independent test performanc before and after<br/>feature selection(Sn, Sp and ACC:%)</b> . . . . .   | <b>36</b> |
| 3.3 | <b>Independent test results on benchmark datasets(Sn, Sp<br/>and ACC:%)</b> . . . . .   | <b>38</b> |
| 3.4 | <b>Cross-validation results on benchmark datasets(Sn, Sp<br/>and ACC:%; TP: true positive, FN: false negative, FP: false<br/>positive, TN: true negative)</b> . . . . . | <b>39</b> |
| 3.5 | <b>4mC site identificaiton in case studies(TP: True Postive;;<br/>FN: False Negative)</b> . . . . .   | <b>41</b> |
| 4.1 | <b>Ranking details of top 12 specific SNP positions (FET:<br/>Fisher's exact test)</b> . . . . .  | <b>53</b> |
| 4.2 | <b>Performance on the independent test dataset (Methy:<br/>Methy-RNA; NPPS: RAM-NPPS)</b> . . . . .   | <b>55</b> |
| 4.3 | <b>Average performance on individual 1226 transcripts<br/>(Methy: Methy-RNA; NPPS: RAM-NPPS)</b> . . . . .  | <b>57</b> |
| 4.4 | <b>Different feature space performance in cross validation<br/>(CPD: Chemical Property with Density; Joint: joint of<br/>conventional features)</b> . . . . .           | <b>58</b> |
| 4.5 | <b>Results for the c-Jun gene case study (Methy: Methy-<br/>RNA; NPPS: RAM-NPPS)</b> . . . . .  | <b>61</b> |

|     |  |    |
|-----|--|----|
| 5.1 | Top ranking expressed genes between two type of issues (NAT: Normal Adjacent Tumor) . . . . .              | 71 |
| 5.2 | Cross-validation performance of top ten genes from different groups (NAT: Normal Adjacent Tumor) . . . . . | 72 |
| 5.3 | Cross-validation performance of top ten genes selected by different DEA methods . . . . .                  | 73 |
| 5.4 | Go function analysis for the top ranking genes (p-value $< 1.0e-04$ and count $\geq 5$ ). . . . .          | 75 |
| 5.5 | Expression boundary of lung cancer biomarkers ( $e$ : FPKM expression level) . . . . .                     | 77 |
| A.1 | The example of 2*2 contingency table . . . . .   | 85 |

# List of Publications

The journal and conference papers published during my PhD study are listed as follows:

## Related to the Thesis :

1. **Z. Zhao**, H. Peng, J. Li et al. Imbalance learning for the prediction of N 6-Methylation sites in mRNAs [J]. BMC genomics, 2018, 19(1), 574.
2. **Z. Zhao**, H. Peng, J. Li et al. Identification of lung cancer gene markers through kernel maximum mean discrepancy and information entropy[J]. BMC Medical Genomics, 2019, 12(8): 1-10.
3. **Z. Zhao**, X. Zhang, J. Li et al. Accurate prediction of DNA N<sup>4</sup>-methylcytosine sites via boost-learning various types of sequence features[J]. BMC genomics, 2020, 21(1), 1-11.

## Others :

4. X. Zhang, **Z. Zhao**, Y. Zheng, and J Li. Prediction of Taxi Destinations Using a Novel Data Embedding Method and Ensemble Learning [J]. IEEE Transactions on Intelligent Transportation Systems, 2019.
5. H. Peng, Y. Zheng, **Z. Zhao**, J. Li et al. Recognition of CRISPR/Cas9 off-target sites through ensemble learning of uneven mis-match distributions

- [J]. *Bioinformatics*, 2018, 34 (17), i757-i765.
6. X. Zhang, Y. Liu, Y. Zheng, **Z. Zhao**, J Li et al. Distinction Between Ships and Icebergs in SAR Images Using Ensemble Loss Trained Convolutional Neural Networks [C]. *Australasian Joint Conference on Artificial Intelligence*. Springer, 2018: 216-223.
  7. Y. Zheng, H. Peng, X. Zhang, **Z. Zhao**, J. Li, et al. Predicting adverse drug reactions of combined medication from heterogeneous pharmacologic databases [J], *BMC Bioinformatics*, 2018, 19(S19).
  8. Y. Zheng, H. Peng, X. Zhang, **Z. Zhao**, J. Li, et al. Old drug repositioning and new drug discovery through similarity learning from drug-target joint feature spaces [J], *BMC bioinformatics*, 2019, 20(23): 605.
  9. Y. Zheng, H. Peng, X. Zhang, **Z. Zhao**, J. Li, et al. DDI-PULearn: a novel positive-unlabeled learning method for large-scale prediction of drug-drug interactions [J], *BMC bioinformatics*, 2019, 20(19): 1-12.



# Abstract

This thesis employs data mining techniques to discover domain knowledge in epigenetic modification and gene expression profile. Computational methods are developed for three research questions, namely, how to accurately predict DNA N<sup>4</sup>-methylcytosine site, how to precisely identify mRNA N<sup>6</sup>-methyladenosine sites, and how to identify lung cancer gene expression profile markers. The motivations of the proposed methods are improving the performance of computational methods via constructing efficient feature space, optimizing machine learning schemes, solving the data imbalance issue, and employing novel statistical analysis approach to provide researchers efficient computational tools.

DNA N<sup>4</sup>-methylcytosine (4mC) is a critical epigenetic modification and plays various roles in the restriction-modification system. The computational methods have been explored to identify 4mC in the DNA sequence in recent years due to the high cost of experimental laboratory detection. However, the state-of-the-art methods have limited performance because of the lack of effective sequence features and the ad hoc choice of learning algorithms. Chapter 3 proposes a new method with novel sequence feature space and machine learning scheme. In sequence encoding, five essential sequence features are integrated into a 292-dimension feature space, representing both global and local sequence characteristics. Then a feature selection scheme is built, where the feature importance score produced from the training process of XGBoost machine is taken as the criterion of feature selection. At last, an SVM-based prediction model is trained with the selected features

and optimized by 10-fold cross-validations. In the result part, the impact of feature selection on model performance is evaluated by an independent test. The proposed method outperforms three state-of-art predictors in both independent test and 10-fold cross-validation. Furthermore, two case studies prove the effectiveness of our method in practical situations.

N<sup>6</sup>-methyladenosine (m<sup>6</sup>A) widely involves in mRNA metabolism and embryogenesis. Multiple computational human mRNA m<sup>6</sup>A site predictors have been developed. However, there are two main drawbacks of the existing methods: first, inadequate learning of the imbalanced training data; second, the sequence text features are not outstanding in representing m<sup>6</sup>A sequence characteristics. Chapter 4 proposes to use the cost-sensitive learning idea to solve the imbalance data issues in the problem. This cost-sensitive approach learns from the entire imbalanced dataset without a random selection of negative samples. In sequence representation, site location, entropy features and specific single nucleotide polymorphism (SNP) positions are taken as new features, which improve the performs significantly. In the comparison with existing predictors, our method achieves better correctness and robustness in both independent tests and case studies. The results suggest that imbalance learning is promising to improve the performance of m<sup>6</sup>A prediction.

The early diagnosis of lung cancer has been a challenging problem in clinical practice for a long time. The identification of differentially expressed genes as a disease marker is a promising solution. Chapter 5 presents a novel approach to identify marker genes and define the boundary of gene expression profile for human lung cancer. By calculating the kernel maximum mean discrepancy, the proposed method evaluates the expression difference between normal, normal adjacent to tumor (NAT) and tumor samples. The expression level boundaries among different groups are defined with the information entropy theory for marker genes. Compared with two conventional methods t-test and fold change, the genes selected by MMD values have better performance under all metrics in 10-fold cross-validation. Furthermore, the GO and KEGG enrichment analysis validate the discovered

marker gene in function pathways. At last, we choose ten most meaningful genes as lung cancer markers and calculate the expression profile boundaries. The proposed method is more accurate than conventional DEA methods in marker gene identification and provides a reliable method for defining the gene expression level boundaries.

